CoCo: Coherence-Enhanced Machine-Generated Text Detection Under Low Resource With Contrastive Learning

Anonymous ACL submission

Abstract

Machine-Generated Text (MGT) detection, a task that discriminates MGT from Human-Written Text (HWT), plays a crucial role in preventing misuse of text generative models, which excel in mimicking human writing style recently. Latest proposed detectors usually take coarse text sequences as input and fine-tune pretrained models with standard cross-entropy loss. However, these methods fail to consider the linguistic structure of texts. Moreover, they lack the ability to handle the low-resource problem which could often happen in practice con-013 sidering the enormous amount of textual data online. In this paper, we present a coherencebased contrastive learning model named CoCo to detect the possible MGT under low-resource 017 scenario. To exploit the linguistic feature, we encode coherence information in form of graph into text representation. To tackle the challenges of low data resource, we employ a contrastive learning framework and propose an im-022 proved contrastive loss for preventing performance degradation brought by simple samples. The experiment results on two public datasets and two self-constructed datasets prove our ap-026 proach outperforms the state-of-art methods significantly. 027

1 Introduction

034

Thriving progress in the field of text generative models (TGMs) (Kenton and Toutanova, 2019; Yang et al., 2019; Liu et al., 2019; See et al., 2019; Keskar et al., 2019; Dathathri et al., 2019; Lewis et al., 2020; Brown et al., 2020; Gao et al., 2021a; Madotto et al., 2021; Ouyang et al., 2022), *e.g.*, ChatGPT¹, enables everyone to produce MGTs massively and rapidly. However, the accessibility to high-quality TGMs is prone to cause misuses, such as fake news generation (Zellers et al., 2019; Brown et al., 2020; Yanagi et al., 2020), product review forging (Adelani et al., 2020), and spamming



Figure 1: Illustration of sentence-level structure difference between HWT and MGT, the MGT is generated by GROVER (Zellers et al., 2019). HWT is more coherent than MGT as the sentences share more same entities with each other.

(Tan et al., 2012), etc. MGTs are hard to distinguish by an untrained human for their human-like writing style (Ippolito et al., 2020) and the excessive amount (Grinberg et al., 2019), which calls for the study of reliable automatic MGT detectors.

Previous works on MGTs detection mainly concentrate on sequence feature representation and classification (Gehrmann et al., 2019; Solaiman et al., 2019; Zellers et al., 2019). Recent studies have shown the good performance of automated detectors in a fine-tuning fashion (Solaiman et al., 2019). Although the fine-tuning based detectors have demonstrated their effectiveness, they still suffer from two issues that limit their conversion to practical use: (1) Existing detectors treat input documents as flat sequences of tokens and use neural encoders or statistical features (*e.g.*, TF-IDF) to represent text as the dense vector for classification. These methods rely much on the token-level dis-

¹https://chat.openai.com

061

081

0

0

0

098 099

100 101

102

103 104

105 106

107

108

109

tribution difference of texts in each class, which ignores high-level linguistic representation of text structure. (2) Compared with the enormous number of online texts, annotated dataset for training MGT detectors is rather low-resource. Constrained by the amount of available annotated data, traditional detectors sustain frustrating accuracy and even collapse during the test stage.

As shown in Fig. 1, MGTs and HWTs exhibit difference in terms of coherence traced by entity consistency. Thus, we propose an entity coherence graph to model the sentence-level structure of texts based on the thoughts of Centering Theory (Grosz and Sidner, 1986), which evaluates text coherence by entity consistency. Entity coherence graph treats entities as nodes and builds edges between entities in the same sentences and same entities among different sentences to reveal the text structure. Instead of treating text as flat sequence, coherence modeling helps to introduce distinguishable linguistic feature at input stage and provides explainable difference between MGTs and HWTs.

To alleviate the low-resource problem in the second issue, inspired by the resurgence of contrastive learning (He et al., 2020; Chen et al., 2020), we utilize proper design of sample pair and contrastive process to learn fine-grained instance-level features under low resource. However, it has been proven that the easiest negative samples are unnecessary and insufficient for model training in contrastive learning (Cai et al., 2020). To circumvent the performance degradation brought by the easy samples, we propose a novel contrastive loss with capability to reweight the effect of negative samples by difficulty score to help model concentrate more on hard samples and ignore the easy samples.

Extensive experiments on multiple datasets (GROVER, GPT-2, GPT-3) and a case study with ChatGPT-generated texts demonstrate the effective-ness of our proposed method.

In summary, our contributions are summarized as follows:

• Coherence Graph Construction: We model the text coherence with entity consistency and sentence interaction while statistically proving its distinctiveness in MGTs detection, and further introduce this linguistic feature at input stage.

• **Improved Contrastive Loss:** We propose a novel contrastive loss in which hard negative

samples are paid more attention for improving detection accuracy of challenging sample.

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

• Outstanding Performance: We achieve stateof-art performance on four MGT datasets in both low-resource and high-resource setting. Experimental results verify the effectiveness of our model.

2 Related Work

Machine-generated Text Detection. Machinegenerated texts, also named deepfake or neural fake texts, are generated by language models to mimic human writing style, making them perplexing for humans to distinguish (Ippolito et al., 2020). Generative models like GROVER (Zellers et al., 2019), GPT-2 (Radford et al., 2019), GPT-3 (Brown et al., 2020) and emerging ChatGPT has been evaluated on the MGT detection task and achieve good results. Bakhtin et al. (2019) train an energy-based model by treating the output of TGMs as negative samples to demonstrate the generalization ability. Deep learning models that incorporate stylometry and external knowledge are also feasible for improving the performance of MGT detectors (Uchendu et al., 2019; Zhong et al., 2020). Our method differs from the previous work by analyzing and modeling text coherence as distinguishable feature and emphasizing the performance improvement under low-resource scenario.

Coherence Modeling. For generative models, coherence is the critical requirement and vital target (Hovy, 1988). Previous works mainly discuss two types of coherence, local coherence (Mellish et al., 1998; Althaus et al., 2004) and global coherence (Mann and Thompson, 1987). Local coherence focus on sentence-to-sentence transitions (Lapata, 2003), while global coherence tries to capture comprehensive structure (Karamanis and Manurung, 2002). Our method strives to represent both local and global coherence with inner- and inter-sentence relations between entity nodes.

Contrastive Learning. Contrastive learning in NLP demonstrates superb performance in learning token-level embeddings (Su et al., 2021) and sentence-level embeddings (Gao et al., 2021b) for natural language understanding. With in-depth study of the mechanism of contrastive learning, the hardness of samples is proved to be crucial in the training stage. Cai et al. (2020) define the dot product between the queries and the negatives



Figure 2: Overview of CoCo. Input document is parsed to construct a coherence graph (3.1), the text and graph are utilized by a supervised contrastive learning framework (3.2), in which coherence encoding module is designed to encode and aggregate to generate coherence-enhanced representation (3.2.3). After that, we employ a MoCo based contrastive learning architecture in which key encodings are stored in a dynamic memory bank (3.2.4) with improved contrastive loss to make final prediction (3.2.5).

in normalized embedding space as hardness and figured out the easiest 95% negatives are insufficient and unnecessary. Song et al. (2022) propose a difficulty measure function based on the distance between classes and apply curriculum learning to the sampling stage. Differently, our method pays more attention to hard negative samples for improving the detection accuracy of challenging samples.

3 Methodology

159

160

161

162

164

165

166

168

169

170

172

173

175

176

177

178

179

180

182

183

The workflow of CoCo mainly contains coherence graph construction, and supervised contrastive learning discriminator, and Fig. 2 illustrates its overall architecture.

3.1 Coherence Graph Construction

In this part, we illustrate how to construct coherence graph to dig out coherence structure of text by modeling sentence interaction.

According to Centering Theory (Grosz and Sidner, 1986), coherence of texts could be modeled by sentence interaction around center entities. To better reflect text structure and avoid semantic overlap, we proposed to construct an undirected graph with entities as nodes. Specifically, we first implement the ELMo-based NER model TagLM (Peters et al., 2017) with the help of NER toolkit AllenNLP² to extract the entities from document. An relation < inter > is constructed between same entities in different sentences and nodes within same sentences are connected by relation < inner >for their natural structure relevance. Formally, the mathematical form of coherence graph's adjacent matrix is defined as follows:

$$\boldsymbol{A}_{ij} = \begin{cases} 1 & rel \langle \texttt{inner} \rangle & v_{i,a} \neq v_{j,b}, a = b \\ 1 & rel \langle \texttt{inter} \rangle & v_{i,a} = v_{j,b}, a \neq b \\ 0 & rel \text{ None} & others \end{cases}$$
191

184

185

186

187

188

189

190

192

193

194

195

196

197

198

199

200

201

202

203

205

206

207

where $v_{i,a}$ represents *i*-th entity in sentence *a*, which is regarded as node in coherence graph.

3.2 Supervised Contrastive Learning

3.2.1 Model Overview

The training process is illustrated in Fig. 2. Each entry in the dataset is document with its coherence graph. The entries in training set are sampled randomly into keys and queries. Two coherence encoder modules (CEM) f_k and f_q , are initialized the same to generate coherence-enhanced representation D_k and D_q for key and query. A dynamic memory bank with the size of all training data is initialized to store all key representation and their annotations for providing enough contrastive pairs in low-resource scenario. In every training step, the newly encoded key graphs update memory bank following First In First Out (FIFO) rule to keep it updated and the training process consistent. A

²https://demo.allennlp.org/named-entity-recognition

novel loss composed of improved contrastive loss 210 and cross-entropy loss ensures the model's ability 211 to achieve instance-level intra-class compactness 212 and inter-class separability while maintaining the 213 class-level distinguishability. A linear discrimina-214 tor takes query representations as input and gener-215 ates prediction results. The pseudocode of training 216 process is shown in Appendix A.8. 217

3.2.2 Positive/Negative Pair Definition

In supervised setting, where we have access to label information, we define two samples with same label as positive pair and that with different labels as negative pair for incorporating label information into training process.

3.2.3 Encoder Design

219

222

225

240

241

In this part, we introduce how to initialize node representation and graph neural network structure which is utilized to integrate coherence information into semantic representation of text by propagating and aggregating information from different granularity with an innovated coherence encoder module (CEM).

Node Representation Initialization. We initialize the representation of entity nodes with powerful pre-trained model RoBERTa for its superior ability to encode contextual information into text representation.

Given an entity e with a span of n tokens, we utilize RoBERTa to map input document x to embeddings h(x). The contextual representation of eis calculated as follows:

$$\boldsymbol{Z}_{v} = \frac{1}{n} \sum_{i=0}^{n} \boldsymbol{h}(\boldsymbol{x})_{e_{i}}, \qquad (1)$$

where e_i is the absolute position where the *i*-th token in *e* lies in the whole document.

Relation-aware GCN. Based on the vanilla Graph Convolutional Networks (Welling and Kipf, 2016), 245 we propose a novel method to assign different weight W_r for inter and inner relation r with 247 Relation-aware GCN. Relation-aware GCN con-248 volute edges of each kind of relation in the coher-249 ence graph separately. The final representation is the sum of GCN outputs from all relations. We use two-layer GCN in the model because more layers will cause an overfitting problem under low resources. We define the relation set as R, and the calculation formula is as follows:

$$H^{(i+1)} = \sum_{r \in R} \hat{A} \text{ReLU}((\hat{A}H^{(i)}W_r^{(i)})W_r^{(i+1)}),$$

$$\hat{A} = \tilde{D}^{-\frac{1}{2}}\tilde{A}\tilde{D}^{-\frac{1}{2}},$$
(2) 256

257

258

259

260

261

262

263

264

265

266

267

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

287

288

289

290

291

292

293

294

where $H^{(i)} \in \mathbb{R}^{N \times d}$ is node representation in *i*-th layer. $\tilde{A} = A + I$, A is the adjacency matrix of the coherence graph, \hat{A} is the normalized Laplacian matrix of \tilde{A} , W_r is the relation transformation matrix for relation r.

Sentence Representation. Afterward, we aggregate updated node representation from last layer of Relation-aware GCN into sentence-level representation to prepare for concatenation with sequence representation from RoBERTa. The aggregation follows the below rule:

$$\boldsymbol{Z}_{s_i} = \frac{1}{M_i} \sum_{j}^{M_i} \sigma(\boldsymbol{W}_s \boldsymbol{H}_{(i,j)} + \boldsymbol{b}_s), \qquad (3)$$

where M_i represents the number of entities in *i*th sentence, $H_{(i,j)}$ represents the embedding of *j*-th entity in *i*-th sentence, W_s is weight matrix and b_s is bias. All the sentence representations within same document are concatenated as sentence matrix Z_s .

Document Representation with Attention LSTM. We design a self-attention mechanism for discovering the sentence-level coherence between one sentence and other sentences, and apply LSTM with the objective to track the coherence in continuous sentences and take the last hidden state of LSTM for aggregated document representation containing comprehensive coherence information. The calculation is described as follows:

$$\boldsymbol{Z}_{c} = \text{LSTM}(\text{softmax}(\gamma \frac{\text{norm}(\boldsymbol{K})\text{norm}(\boldsymbol{Q})^{T}}{\sqrt{d_{Z}}})\boldsymbol{V}), \quad (4)$$

where K, Q, V are linear transformations of Z_s with matrix W_k, W_q, W_v, d_Z is the dimension of representation Z_s , and γ is a hypergammarparameter for scaling.

Finally, we concatenate Z_c and the sequence representation h([CLS]) from the RoBERTa's last layer to generate document coherence-enhanced representation D.

3.2.4 Dynamic Memory Bank

The dynamic memory bank is created to store as much as key encoding D_k to form adequate positive and negative pairs within a batch. The dynamic



Figure 3: Illustration of coherence encoder module (CEM) which encodes and fuses the coherence graph and text sequence to generate coherence-enhanced representation of document.

memory bank is maintained as a queue so that the newly encoded keys could replace the outdated ones, which keeps the consistency between the key encoding and current training step.

3.2.5 Loss Function

297

298

307

310

311

314

315

318

Following the definition of positive pairs and negative pairs above, traditional supervised contrastive loss (Gunel et al., 2020) treats all positive pairs and negative pairs equally.

However, with recognition that not all negatives are created equal (Cai et al., 2020), our goal is to emphasize the informative samples for helping the model to differentiate difficult samples. Thus, we propose an improved contrastive loss which dynamically adjusts the weight of negative pair similarity according to the hardness of negative samples. To be specific, the hard negative samples should be assigned larger weight for stimulating the model to better pull same class together and push different class away. The improved contrastive loss is defined as:

$$\mathcal{L}_{\text{ICL}} = \sum_{j=1}^{M} \mathbf{1}_{y_i = y_j} \log \frac{S_{ij}}{\sum_{p \in \mathcal{P}(i)} S_{ip} + \sum_{n \in \mathcal{N}(i)} rf_{in}S_{in}},$$
$$rf_{ij} = \beta \frac{\mathbf{D}_q^i \mathbf{D}_n^n}{\operatorname{avg}(\mathbf{D}_q^i \mathbf{D}_k^{n+1})},$$
$$S_{ij} = \exp(\mathbf{D}_q^i \mathbf{D}_k^j / \tau),$$
(5)

where $\mathcal{P}(i)$ is the positive set in which data has the same label with q_i and $\mathcal{N}(i)$ is the negative set in which data has different label from q_i .

> Apart from instance-level learning mechanism, a linear classifier combined with cross entropy loss \mathcal{L}_{CE} is employed to provide the model with class

level separation ability. \mathcal{L}_{CE} is calculated by

$$\mathcal{L}_{CE} = \frac{1}{N} \sum_{i=1}^{N} -[y_i log(p_i) + (1 - y_i) log(1 - p_i)],$$
(6)

where p_i is the prediction probability distribution of *i*-th sample. The final loss $\mathcal{L}_{\text{total}}$ is a weighted average of \mathcal{L}_{ICL} and \mathcal{L}_{CE} as:

$$\mathcal{L}_{\text{total}} = \alpha \mathcal{L}_{\text{ICL}} + (1 - \alpha) \mathcal{L}_{\text{CE}}, \qquad (7)$$

where the hyperparameter α adjusts the relative balance between instance compactness and class separability.

3.2.6 Momentum Update

ł

The parameters of query encoder f_q and the classifier can be updated by gradient back-propagated from $\mathcal{L}_{\text{total}}$. We denote the parameters of f_q as θ_q , the parameters of f_k as θ_k , The key encoder f_k 's parameters are updated by momentum update mechanism:

$$\theta_k \leftarrow \beta \theta_k + (1 - \beta) \theta_q,$$
(8)

where the hyperparameter β is momentum coefficient.

4 Experiments

4.1 Datasets

We evaluate our model on the following datasets:

GROVER Dataset is a News-Style dataset provided by (Zellers et al., 2019), in which 348 HWTs are collected from RealNews, a large corpus of news articles from Common Crawl, 350 and MGTs are generated by Grover-Mega, a 351 transformer-based news generator with parameters sized 1.5B.

326

327

328

329

330

331

333

334

335

336

337

338

339

341

344

424

425

426

427

428

429

430

431

432

433

434

435

436

437

389

390

391

Dataset		Train	Valid	Test
CDOVED	HWT	5,000	2,000	8,000
GROVER	MGT	5,000	1,000	4,000
CDT 2	HWT	25,000	5,000	5,000
GP1-2	MGT	25,000	5,000	5,000
CDT 2 unmixed	HWT	5,455	1,000	1,000
GF I-5 ullilized	MGT	5,455	1,000	1,000
GPT-3 mixed	HWT	5,033	1,000	1,000
	MGT	5,033	1,000	1,000

Table 1: Basic statistics of datasets.

- **GPT-2 Dataset** is a Webtext-style dataset provided by OpenAI³ with HWTs adopted from WebText and MGTs produced by GPT-2 XLM-1542M.
- GPT-3 Dataset is a News-Style open source dataset constructed by us based on the text-davinci-003⁴ model of OpenAI, which is the most capable GPT-3 model so far and can generate longer texts (maximum 4,000 tokens). The GPT-3 model refer to various latest newspapers (Dec. 2022 Present) whose full texts act as the HWTs part and generate news by imitation. We use two subsets: mixed- and unmixed-provenances. The details of this dataset are explained in Appendix A.1.

The statistics of datasets is summarized in Table 1 and the implementation details are in Appendix A.3.

4.2 Comparison Models

355

357

358

363

364

369

370

371

372

373

374

375

378

386

387

We compare CoCo to the state-of-art detection methods to reveal the effectiveness. The baselines are introduced as follows.

- **GPT-2** (Radford et al., 2019), a strong pretrained language model based on the decoder architecture of transformer. We use GPT-2 small with parameters sized 124M for the fairness of comparison.
- **RoBERTa** (Liu et al., 2019), a powerful transformers-based bidirectional language model with robust performance on downstream tasks. We use RoBERTa-base sized 110M in the experiment.
- XLNet (Yang et al., 2019), a language model which is superb in understanding long docu-

ments. We exploit XLNet-base whose scale is 110M.

- **CE+SCL** (Gunel et al., 2020), a state-of-theart supervised contrastive learning method in various downstream task. We train the detector with Cross-Entropy loss (CE) and supervised contrastive loss (SCL) calculated within a mini-batch.
- **DualCL** (Chen et al., 2022), a contrastive learning method with the addition of label representations for data augmentation.

4.3 Performance Comparison

As shown in Table 2, COCO surpasses state-of-theart methods in MGT detection task 4.25%, 5.98% (limited dataset), and 3.31%, 3.58% (full dataset) in average in terms of accuracy and F1-score, respectively. The result indicates the rationality of coherence graph designing and the effectiveness of encoding coherence information into document representation, which further proves the coherence difference between MGT and HWT. Moreover, it should be noticed that CoCo gets less affected by randomness, which illustrates that the coherence graph we construct is a robust feature that helps stabilize the model performance. Meanwhile, CoCo outperforms CE+SCL and DualCL regardless the size of training set, which suggests the success of improved contrastive loss to solve the performance degradation problem brought by simple negative samples.

We also find GROVER Dataset is the hardest to detect. It is because GROVER generator is trained in an adversarial fashion with the objective to deceive the verifier, which endows the generator with deceptive nature. To our surprise, GPT-3 dataset is overly simple to all detectors. We conduct extensive experiments on different self-constructed and published GPT-3 dataset generated by a series prompts, which also validate this thundering conclusion. The experiment details and results are in Appendix A.2.

Due to the secrecy of GPT-3 implementation details and the lack of interpretability of LLMs, we can only make several inferences for this phenomenon up to the best of our knowledge. First, since GPT-3 is fine-tuned with RLHF (Ouyang et al., 2022), we suppose that instructions like "intimate a piece of news" is rare in prompt dataset used in GPT-3 fine-tuning, which causes its inadequate ability to imitate news. In fact, we find news

³https://github.com/openai/gpt-2-output-dataset

⁴https://beta.openai.com/docs/models/gpt-3

Dataset		GROVER				GPT-2			
Size	Limited Da	taset (10%)	Full D	Dataset	Limited Dataset (10%)		Full Dataset		
Metric	ACC	F1	ACC	F1	ACC	F1	ACC	F1	
GPT2	0.6401 ± 0.0136	0.4289 ± 0.0413	0.8274 ± 0.0091	0.8003 ± 0.0141	0.8575 ± 0.0041	0.8406 ± 0.0070	0.8913 ± 0.0066	0.8839 ± 0.0078	
XLNet	0.6906 ± 0.0321	0.5193 ± 0.0365	0.8156 ± 0.0079	0.7493 ± 0.0073	0.8837 ± 0.0031	0.8732 ± 0.0041	0.9091 ± 0.0091	0.9027 ± 0.0111	
RoBERTa	0.7730 ± 0.0121	0.6370 ± 0.0186	0.8772 ± 0.0029	0.8171 ± 0.0048	0.9244 ± 0.0041	0.9214 ± 0.0045	0.9402 ± 0.0039	0.9384 ± 0.0044	
DualCL	0.6926 ± 0.0634	0.5397 ± 0.0971	0.7574 ± 0.0855	$\textbf{0.6388} \pm \textbf{0.130}$	0.7874 ± 0.1210	0.6950 ± 0.0944	0.8023 ± 0.1120	0.8046 ± 0.1530	
CE+SCL	0.7777 ± 0.0134	0.6447 ± 0.0286	0.8782 ± 0.0044	0.8202 ± 0.0057	0.9255 ± 0.0030	0.9241 ± 0.0016	0.9408 ± 0.0006	0.9390 ± 0.0009	
CoCo	$\textbf{0.7808} \pm \textbf{0.0044}$	$\textbf{0.6543} \pm \textbf{0.0106}$	$\textbf{0.8826} \pm \textbf{0.0018}$	$\textbf{0.8265} \pm \textbf{0.0036}$	$\textbf{0.9271} \pm \textbf{0.0019}$	$\textbf{0.9254} \pm \textbf{0.0018}$	$\textbf{0.9457} \pm \textbf{0.0004}$	$\textbf{0.9452} \pm \textbf{0.0004}$	
Dataset		GPT-3 U	Unmixed		GPT-3 Mixed				
Size	Limited Da	taset (10%)	Full D	Dataset	Limited Da	Limited Dataset (10%) Full		ataset	
Metric	ACC	F1	ACC	F1	ACC	F1	ACC	F1	
GPT2	0.9631 ± 0.0037	0.9666 ± 0.0044	0.9917 ± 0.0056	0.9905 ± 0.0042	0.9623 ± 0.0039	0.9636 ± 0.0046	0.9910 ± 0.0046	0.9910 ± 0.0033	
XLNet	0.9252 ± 0.0052	0.9241 ± 0.0054	0.9620 ± 0.0043	0.9634 ± 0.0068	0.9149 ± 0.0044	0.9152 ± 0.0059	0.9513 ± 0.0052	0.9505 ± 0.0039	
RoBERTa	0.9916 ± 0.0043	0.9899 ± 0.0049	0.9928 ± 0.0035	0.9913 ± 0.0040	0.9892 ± 0.0022	0.9893 ± 0.0025	0.9923 ± 0.0017	0.9901 ± 0.0024	
CE+SCL	0.9918 ± 0.0039	0.9901 ± 0.0058	0.9944 ± 0.0023	0.9943 ± 0.0031	0.9903 ± 0.0020	0.9898 ± 0.0023	0.9932 ± 0.0017	0.9905 ± 0.0038	
CoCo	$\textbf{0.9932} \pm \textbf{0.0042}$	$\textbf{0.9913} \pm \textbf{0.0033}$	$\textbf{0.9972} \pm \textbf{0.0015}$	$\textbf{0.9957} \pm \textbf{0.0020}$	$\textbf{0.9913} \pm \textbf{0.0034}$	$\textbf{0.9911} \pm \textbf{0.0038}$	$\textbf{0.9932} \pm \textbf{0.0019}$	$\textbf{0.9937} \pm \textbf{0.0028}$	

Table 2: Results of the model comparison. It should be noticed that DualCL is easily affected by random seed, which may be caused by its weakness in understanding long texts. We do not present the experiment results for DualCL on GPT-3 dataset because the documents in GPT-3 dataset is so long that DualCL completely fails.

generated by GPT-3 is more like a short version 438 439 of reference news (far shorter than the max output length) instead of the imitation. Second, the HWTs 440 we use are latest news which are not included in 441 GPT-3's training data. The lack of knowledge back-442 443 ground about reference news might limit GPT-3's associative ability. Third, we do not fine-tune GPT-3 444 with HWTs we collect while model used in GPT-2 445 dataset did. The unfamiliarity of GPT-3 with the 446 texts it is required to imitate could impair its ability 447 to rewrite news in a similar style. We will further 448 investigate our hypotheses and provide possible 449 deep insights into this counterintuitive but very in-450 451 teresting phenomenon, which may also exist in the GPT-4 model, in future works. 452

4.4 Ablation Study

453

454

455

456

457

To illustrate the necessity of some components of CoCo, we conduct several ablation experiments on limited GROVER dataset and the results are shown in Table 3. We introduce the ablation model structure below:

Model	ACC	F1
CoCo (Plain)	0.7697	0.6428
CoCo (Sentence nodes)	0.7733	0.6379
CoCo (Coherence)	0.7777	0.6463
CoCo (Coherence + LSTM)	0.7787	0.6471
CoCo (Coherence + LSTM + SCL)	0.7827	0.6609
СоСо	0.7843	0.6684

Table 3: Results of ablation study.

CoCo (Plain) removes graph information and encodes only by RoBERTa parts. Moreover, the model removes contrastive learning and uses CE loss.

CoCo (Sentence Nodes) treats sentences as nodes instead of entities and establish edges between sentences which share same entities. Node representation is initialized by RoBERTa embedding and mean-pooling operation. Document representation is obtained by one CEM discarding sentence representation and attention LSTM part in section 3.2.3. Document representation is calculated by meanpooling operation on sentence node representations. A linear classification head with cross-entropy loss is used for detection.

COCO (Coherence) incorporates coherence graph into the representation of document and deploys sentence representation part in section 3.2.3. The rest are the same with COCO (Sentence Nodes).

COCO (Coherence + LSTM) uses attention LSTM for document-level aggregation and the rest is the same as COCO (Coherence).

COCO (Coherence + LSTM + SCL) utilizes the contrastive learning framework but the loss function is traditional supervised contrastive loss instead of the improved contrastive loss.

As shown in Table 3, coherence information and the contrastive learning framework greatly contribute to the development of model performance, especially in F1-Score. Replacing entity nodes in

485

486

487

458

459

460

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

549

550

551

552

553

554

555

518

coherence graph with sentences impairs the detector, which could be caused by semantic overlap between graph representation and text sequence representation. The attention LSTM also plays an important role in preserving coherence information during sentence aggregation. Lastly, the results also shows the advantage of improved contrastive loss over standard supervised contrastive loss. We also investigate the effect of hyperparameters on COCO, the results are shown in Appendix A.4.

Discussion: Utility of Coherence Graph 4.5

Case Study 4.5.1

488

489

490

491

492

493

494

495

496

497

498

499

507

509

510

511

512

513

514

515

In this subsection, we conduct a case study with HWT and MGT produced by sensational ChatGPT with the same metadata. As illustrated in Fig. 4, we parse two news as coherence graphs. And we observe that although ChatGPT expresses fluently, it is not coherent from the perspective of coherence graph. Hence, COCO utilizes the distinctive coherence feature and makes correct predictions while RoBERTa fails. This reflects even the most popular and advanced language model could suffer from weak coherence and be detected by CoCo.



Figure 4: An illustration for case study of our method. Entities in documents are colored green. The blue solid box indicates the sentence. The orange dashed lines are inner edges and green dashed lines are inter edges. Numbers in red indicate the probability of predicted label.

4.5.2 Static Analysis for Coherence Graph

We apply static geometric features analysis on coherence graph we construct to provide an in-depth view on linguistic explanation about graph structure. In the following discussion, we take the

dataset of GROVER into the analysis. Some basic metrics of data and the corresponding graph are shown in Table 9.

Metric	HWT	MGT
Sample Num.	4994	4991
Avg. Num. of Token	463.2	456.0
Avg. Num. of Vertex	43.60	32.37
Avg. Num. of Edge	107.4	65.44

Table 4: Basic metrics of texts and corresponding coherence graphs.

Degree Distribution of coherence graph semantically measures the co-occurrence and TF-IDF feature of keywords, showing global coherence because high-degree nodes devote to the main topic and low-degree nodes are the extension. The degree of the graph representation of HWTs is 2.980, which is 15.0% larger than MGTs (2.591), and the distribution of HWTs has a longer tail than MGTs. Furthermore, we prove that degree distribution can robustly detect MGTs and HWTs when impacted by style and genre differences. More details are discussed in the Appendix A.5.

5 Conclusion

In this paper, we propose CoCo, a coherenceenhanced contrastive learning model for MGT detection. We construct a novel coherence graph from document and implement a MoCo-based contrastive learning framework to improve model performance in low-resource setting. An innovative encoder composed of relation-aware GCN and attention LSTM is designed to learn the coherence representation from coherence graph which is further incorparated with sequence representation of document. To alleviate the effect of unnecessary easy samples, we propose an improved contrastive learning loss to force the model to pay more attention to hard negative samples. We evaluate our method on MGT datasets generated by GROVER, GPT-2, and GPT-3, respectively, in both lowresource and high-resource settings. CoCo outperforms Transformer-based methods and contrastivelearning-based methods on all datasets and both settings.

Acknowledgements

We thank the six anonymous reviewers and the area chair for their helpful comments and feedback, which aided us in greatly improving the paper.

Limitations

556

586

587

589

593

599

In this work, we step forward to better distinguish-557 ing MGTs under the low-resource setting. How-558 ever, several limitations still exist for the broader 559 applications of this detector. Firstly, MGTs are easier to generate and collect than HWTs, which may cause an imbalanced label distribution in the dataset. And CoCo literally corrupts in extremely imbalanced data distribution condition, as shown in A.6. Future work could build upon the con-565 trastive learning method of CoCo with innovation on sampling strategy for harsh low-resource and imbalanced data settings. Secondly, our method ar-569 tificially generates a coherence graph for every entry, which is not efficient for larger datasets. What's 570 more, short text, codes, and mathematical proofs, which are hard to generate coherence graphs, are also limitedly detected by CoCo. More distinctive and easy-to-calculate features are worth exploring 574 for generating distinguishable representations for texts with efficiency while better understanding the essence of TGMs. Thirdly, with instruct-based 577 generation and human-in-loop fine-tuning models prevailing, the strategy and defect of TGMs change 579 slightly but constantly. The entity relation with the same semantic granularity and concretization 581 in this paper would not be enough to detect the high-quality content by TGMs in the future. More 583 generative and adaptive detection models should be considered. 585

Ethical Considerations

We provide insight into the potential weakness of TGMs and publish GPT-3 news dataset. We understand that the discovery of our work can be viciously used to confront detectors. And we understand that malicious users can copy the contents of our GPT-3 news dataset to disguise real news and publish them. However, with the purpose of calling for attention to detecting and controlling possible misuse of TGMs, we believe our work will inspire the advance of the stronger detector of MGTs and prevent all potential negative uses of language models.

Our work complies with sharing & publication policy of OpenAI⁵ and all data we collect is in public domain and licensed for research purposes.

References

David Ifeoluwa Adelani, Haotian Mai, Fuming Fang, Huy H Nguyen, Junichi Yamagishi, and Isao Echizen. 2020. Generating sentiment-preserving fake online reviews using neural language models and their human-and machine-based detection. In *International Conference on Advanced Information Networking and Applications*, pages 1341–1354. Springer. 602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

- Ernst Althaus, Nikiforos Karamanis, and Alexander Koller. 2004. Computing locally coherent discourses.
 In Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04), pages 399–406.
- Anton Bakhtin, Sam Gross, Myle Ott, Yuntian Deng, Marc'Aurelio Ranzato, and Arthur Szlam. 2019. Real or fake? learning to discriminate machine from human generated text. *arXiv preprint arXiv:1906.03351*.
- Roi Blanco and Christina Lioma. 2011. Graph-based term weighting for information retrieval. *Information Retrieval*, 15:54–92.
- Stefan Bordag, Gerhard Heyer, and Uwe Quasthoff. 2003. Small worlds of concepts and other principles of semantic search. In *International Workshop on Innovative Internet Community Systems*, pages 10–19. Springer.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Tiffany Tianhui Cai, Jonathan Frankle, David J Schwab, and Ari S Morcos. 2020. Are all negatives created equal in contrastive instance discrimination? *arXiv preprint arXiv:2010.06682*.
- Ramon Ferrer Cancho and Ricard V. Solé. 2001. Two regimes in the frequency of words and the origins of complex lexicons: Zipf's law revisited*. *Journal of Quantitative Linguistics*, 8:165 – 173.
- Ramon Ferrer Cancho, Ricard V Solé, and Reinhard Köhler. 2004. Patterns in syntactic dependency networks. *Physical Review E*, 69(5):051915.
- Qianben Chen, Richong Zhang, Yaowei Zheng, and Yongyi Mao. 2022. Dual contrastive learning: Text classification via label-aware data augmentation. *arXiv preprint arXiv:2201.08702*.
- Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. 2020. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*.
- Monojit Choudhury, Markose Thomas, Animesh Mukherjee, Anupam Basu, and Niloy Ganguly. 2007. How difficult is it to develop a perfect spell-checker? a cross-linguistic analysis through complex network

⁵https://openai.com/api/policies/sharing-publication/

- 657 664 670 671 672 673 679 681 685 694 695
- 700 701 702

- 706 707

- 711

approach. In Proceedings of the Second Workshop on TextGraphs: Graph-Based Algorithms for Natural Language Processing, pages 81–88.

- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2019. Plug and play language models: A simple approach to controlled text generation. arXiv preprint arXiv:1912.02164.
- Sergey N Dorogovtsev and José Fernando F Mendes. 2001. Language as an evolving word web. Proceedings of the Royal Society of London. Series B: Biological Sciences, 268(1485):2603-2606.
- Ramon Ferrer Cancho, Andrea Capocci, and Guido Caldarelli. 2007. Spectral methods cluster words of the same class in a syntactic dependency network. International journal of bifurcation and chaos, 17(7):2453-2463.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021a. Making pre-trained language models better few-shot learners. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3816-3830.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021b. Simcse: Simple contrastive learning of sentence embeddings. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 6894-6910.
- Sebastian Gehrmann, Hendrik Strobelt, and Alexander M Rush. 2019. Gltr: Statistical detection and visualization of generated text. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, pages 111-116.
- Nir Grinberg, Kenneth Joseph, Lisa Friedland, Briony Swire-Thompson, and David Lazer. 2019. Fake news on twitter during the 2016 us presidential election. Science, 363(6425):374-378.
- Barbara J Grosz and Candace L Sidner. 1986. Attention, intentions, and the structure of discourse. Computa*tional linguistics*, 12(3):175–204.
- Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. 2020. Supervised contrastive learning for pretrained language model fine-tuning. arXiv preprint arXiv:2011.01403.
- Geehan Sabah Hassan, Asma Khazaal Abdulsahib, and Siti Sakira Kamaruddin. 2017. Graph-based text representation: a survey of current approaches. Research Journal of Applied Sciences, Engineering and Technology, 14(9):334-340.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 9729–9738.

Xiaochen Hou, Peng Qi, Guangtao Wang, Rex Ying, Jing Huang, Xiaodong He, and Bowen Zhou. 2021. Graph ensemble learning over multiple dependency trees for aspect-level sentiment classification. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2884-2894.

712

713

714

716

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

- Eduard H Hovy. 1988. Planning coherent multisentential text. In Proceedings of the 26th annual meeting on Association for Computational Linguistics, pages 163-169.
- Binxuan Huang and Kathleen M Carley. 2019. Syntaxaware aspect level sentiment classification with graph attention networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 5469–5477.
- Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. 2020. Automatic detection of generated text is easiest when humans are fooled. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1808-1822.
- Nikiforos Karamanis and Hisar Maruli Manurung. 2002. Stochastic text structuring using the principle of continuity. In Proceedings of the International Natural Language Generation Conference, pages 81–88.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of NAACL-HLT, pages 4171–4186.
- Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. arXiv preprint arXiv:1909.05858.
- Zornitsa Kozareva, Ellen Riloff, and Eduard Hovy. 2008. Semantic class learning from the web with hyponym pattern linkage graphs. In Proceedings of ACL-08: HLT, pages 1048–1056.
- Mirella Lapata. 2003. Probabilistic text structuring: Experiments with sentence ordering. In ACL, volume 3, pages 545-552. Citeseer.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7871-7880.
- Xien Liu, Xinxin You, Xiao Zhang, Ji Wu, and Ping Lv. 2020. Tensor graph convolutional networks for text classification. In Proceedings of the AAAI conference on artificial intelligence, volume 34, pages 8409-8416.

Ilya Loshchilov and Frank Hutter. 2017.

based learning for dialogue systems.

and mining 2015, pages 1473-1479.

Sciences Institute Los Angeles.

processing, pages 404–411.

arXiv:1711.05101.

abs/2110.08118.

ACL Anthology.

E, 65(6):065102.

arXiv:2203.02155.

blog, 1(8):9.

ing.

2022.

pled weight decay regularization. arXiv preprint

Andrea Madotto, Zhaojiang Lin, Genta Indra Winata,

Fragkiskos D Malliaros and Konstantinos Skianis. 2015.

Graph-based term weighting for text categorization. In Proceedings of the 2015 IEEE/ACM international

conference on advances in social networks analysis

William C Mann and Sandra A Thompson. 1987. Rhetorical structure theory: A theory of text organiza-

Chris Mellish, Alistair Knott, Jon Oberlander, and

Mick O'Donnell. 1998. Experiments using stochas-

tic search for text planning. In Proceedings of the

9th International General Workshop, pages 98–107.

Rada Mihalcea and Paul Tarau. 2004. Textrank: Bring-

ing order into text. In Proceedings of the 2004 con-

ference on empirical methods in natural language

Marvin Minsky. 1982. Semantic information process-

Adilson E Motter, Alessandro PS De Moura, Ying-

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang,

Sandhini Agarwal, Katarina Slama, Alex Ray, et al.

structions with human feedback. arXiv preprint

Matthew E Peters, Waleed Ammar, Chandra Bhaga-

vatula, and Russell Power. 2017. Semi-supervised

sequence tagging with bidirectional language models.

In Proceedings of the 55th Annual Meeting of the

Association for Computational Linguistics (Volume

Alec Radford, Jeffrey Wu, Rewon Child, David Luan,

Dario Amodei, Ilya Sutskever, et al. 2019. Language

models are unsupervised multitask learners. OpenAI

1: Long Papers), pages 1756–1765.

Training language models to follow in-

Cheng Lai, and Partha Dasgupta. 2002. Topology of the conceptual network of language. Physical Review

tion. University of Southern California, Information

and Pascale Fung. 2021. Few-shot bot: Prompt-

- 781
- 785
- 789
- 790 791 792

- 795

798

800

802

- 810
- 811
- 812 813

814

815

817

819

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-Abigail See, Aneesh Pappu, Rohun Saxena, Akhila dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Yerukola, and Christopher D Manning. 2019. Do Luke Zettlemoyer, and Veselin Stoyanov. 2019. massively pretrained language models make better Roberta: A robustly optimized bert pretraining apstorytellers? In Proceedings of the 23rd Conferproach. arXiv preprint arXiv:1907.11692. ence on Computational Natural Language Learning (CoNLL), pages 843-861.

Decou-

ArXiv.

Mariano Sigman and Guillermo A Cecchi. 2002. Global organization of the wordnet lexicon. Proceedings of the National Academy of Sciences, 99(3):1742–1747. 821

822

823

824

825

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

844

845

846

847

848

849

850

851

852

853

854

855

856

857

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

- Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, et al. 2019. Release strategies and the social impacts of language models. arXiv preprint arXiv:1908.09203.
- Xiaohui Song, Longtao Huang, Hui Xue, and Songlin Hu. 2022. Supervised prototypical contrastive learning for emotion recognition in conversation. arXiv preprint arXiv:2210.08713.
- Mark Steyvers and Joshua B Tenenbaum. 2005. The large-scale structure of semantic networks: Statistical analyses and a model of semantic growth. Cognitive science, 29(1):41-78.
- Yixuan Su, Fangyu Liu, Zaiqiao Meng, Tian Lan, Lei Shu, Ehsan Shareghi, and Nigel Collier. 2021. Tacl: Improving bert pre-training with token-aware contrastive learning. arXiv preprint arXiv:2111.04198.
- Enhua Tan, Lei Guo, Songqing Chen, Xiaodong Zhang, and Yihong Zhao. 2012. Spammer behavior analysis and detection in user generated content on social networks. In 2012 IEEE 32nd International Conference on Distributed Computing Systems, pages 305-314. IEEE.
- Peter D Turney. 2002. Learning to extract keyphrases from text. arXiv preprint cs/0212013.
- Adaku Uchendu, Jeffrey Cao, Qiaozhi Wang, Bo Luo, and Dongwon Lee. 2019. Characterizing man-made vs. machine-made chatbot dialogs. In TTO.
- Adaku Uchendu, Zeyu Ma, Thai Le, Rui Zhang, and Dongwon Lee. 2021. TURINGBENCH: A benchmark environment for Turing test in the age of neural text generation. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 2001–2016, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Feng Wang and Huaping Liu. 2021. Understanding the behaviour of contrastive loss. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 2495-2504.
- Max Welling and Thomas N Kipf. 2016. Semisupervised classification with graph convolutional networks. In J. International Conference on Learning Representations (ICLR 2017).

- 874 875
- 876

890

892

896

897

900

901

902

903

904

905

Dominic Widdows and Beate Dorow. 2002. A graph model for unsupervised lexical acquisition. In COL-ING 2002: The 19th International Conference on Computational Linguistics.

- Yuta Yanagi, Ryohei Orihara, Yuichi Sei, Yasuyuki Tahara, and Akihiko Ohsuga. 2020. Fake news detection with generated comments for news articles. In 2020 IEEE 24th International Conference on Intelligent Engineering Systems (INES), pages 85–90. IEEE.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. Advances in neural information processing systems, 32.
 - Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Graph convolutional networks for text classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 7370–7377.
 - Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. Advances in neural information processing systems, 32.
 - Wanjun Zhong, Duyu Tang, Zenan Xu, Ruize Wang, Nan Duan, Ming Zhou, Jiahai Wang, and Jian Yin.
 2020. Neural deepfake detection with factual structure of text. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2461–2470.

A Appendix

A.1 Details of GPT-3 Dataset

GPT-3 Dataset for CoCo is our latest dataset for 906 the MGT detection task. There are two subsets 907 in the self-made dataset for easy analysis of the 908 impact of provenance and writing styles: unmixed-909 and mixed provinces. We use the text-davinci-003 910 model of OpenAI to generate MGT examples. The 911 maximum length of HWTs is 1024 tokens, and the 912 target generation length is set as 1024 tokens. Here 913 is an example of the MGT data. 914

"title": "On Eve of World Cup, FIFA Chief Says, 'Don't Criticize Qatar; Criticize Me.'" "text": "DOHA, Qatar. The president of world soccer's governing body on Saturday sought to blunt mounting concerns about the World Cup in Qatar with a strident defense of both the host country's reputation and FIFA's authority over its showpiece championship..... Citing statistics, history and even childhood to bolster his case, he at one point likened his own experience as a redheaded child of immigrants to Switzerland to the assimilation problems of gays in the Middle East, and defended the laws, customs and honor of the host country.", "authors": ["Tariq Panja"], "publish_date": "2022-11-19 00:00:00", "source": "The New York Times", "url": "https://www.nytimes.com/2022/11/19/sports/ soccer/world-cup-gianni-infantino-fifa.html"

And the following data shows the corresponding MGT in the dataset.

"title": "On Eve of World Cup, FIFA Chief Says, 'Don't Criticize Qatar; Criticize Me.' "text": "The 2022 FIFA World Cup in Qatar is fast approaching, and its organizing committee's president, Gianni Infantino, is speaking out about the lingering criticism of the country hosting the event. he said. "It is a once-in-a-lifetime opportunity for the region to show the world its values and aspirations, and it is vital that this event is seen as a celebration of football and a celebration of the region."" "authors": "machine", "source": "The New York Times", "matched_hwt_id": 202, "label": "machine""

A.1.1 Human Written Texts

Unmixed Subset. The HWTs of the unmixed subset are all from The New York Times⁶ to exclude the impact of writing style. The time span of our data is Nov 1, 2022 - Dec 25, 2022, making sure that no pre-trained model has learned them. We develop the crawler based on news-crawler⁷.

Mixed Subset. The HWTs of the mixed subset come from various sources, listed as Table 5. The time span of the data is Jan 1, 2022 - Jan 7, 2023. We develop the crawler based on Newspaper3k⁸.

The dataset is specifically designed for MGTs detection and improving generation models. The contents of dataset are obtained from official news websites and the names of indicidual people are not mentioned maliciously. And we strongly reject using our dataset to create offensive content or peek at private information.

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

916

917

⁶https://www.nytimes.com/

⁷https://github.com/LuChang-CS/news-crawler

⁸https://github.com/codelucas/newspaper

Name	Website	Generate a news passage.
Kotaku	https://kotaku.com	The news is written by {Auth
The Daily World	https://www.thedailyworld.com	Title: Lionel Messi isn't e
CNN	https://edition.cnn.com	with PSG until early Januar
BBC	https://www.bbc.com	success
NBC News	https://www.nbcnews.com	Keywords: exploring, mounta
Reuters	https://www.reuters.com	Malavath, Kavya Manyapu, NA Project Shakthi girls' A
Huffpost	https://www.huffpost.com	India, virgin peak, climbing,
Pando	http://pandodaily.com	safety, motivation, empower
Yahoo	https://news.yahoo.com	gender gap, Mount Aconcagua,
Sun Times	https://chicago.suntimes.com/news	Entities: CNN, Poorna Malava
Sfgate	https://www.sfgate.com	India. Mount Aconcagua. Sou
New Republic	https://newrepublic.com	Bank.
Time	https://time.com	Examples: designing space
Pcmag	http://www.pcmag.com	ever woman to summit Mount
CNBC	https://www.cnbc.com/world/	education difficulties of
News	https://www.news.com.au/	peak, experiences of altitude
The Atlantic	https://www.theatlantic.com/latest/	of Project Shakthi, India's
		Act, sponsorship for under

Table 5: Data sources for the mixed subset.

A.1.2 **Machine Generated Texts**

As the GPT-3 and ChatGPT model need prompts to generate, we write hints for the generation models to generate texts that meet our news-style long text generation. The hints format is as follows, and the content is related to HWTs.

Write a news more than 1000 words. The news is written by {Authors} from {Source} in {date}. Title is {title}.

A.2 GPT-3 Dataset Generated by Different **Prompts and Experiment Results**

To further validate the conclusion that GPT-3 generated texts are easier to detect, we utilize CNN news as reference and design different prompts for GPT-3 generation. The principle is to provide as more information as possible to GPT-3 for alleviating the possible gap in semantics and in length.

Keywords as Prompt (KP). We extract the keywords and entities with GPT-3.5-turbo and provide examples in original news to form the prompt for generation. The prompt format is as follows.

Example prompt for generation.

"role": "syste	m", "con	tent"	: "Extra	act	all
the keywords, e	entities,	and	examples	in	the
following passag	ge:"				
"role": "user",	"content"	': {te	ext}		

Example prompt for generation.

hors} from {Source} xpected to be back ry after World Cup ins, space, Poorna SA, Mount Everest, education. Ladakh altitude sickness erment, education, sponsorship. ath, Kavya Manyapu, t Shakthi, Ladakh uth America, World suits, youngest Everest, climbed a oney to fund girls' climbing a virgin e sickness, purpose Right to Education rprivileged school children, scaling Mount Aconcagua, expanding sponsorship globally.

The target length for generation is 731 tokens Add as much details and examples as you can. News:

Summary as Prompt (SP). We employ GPT-3.5turbo to summarize the original texts. The compression ratio is set to [0.3, 1.0], which means the summary is required to be longer than 0.3 of the length of original text and shorter than whole original text. The generated summary is used as prompt and the format is as follows:

Generate following a news based on the abstract: Paris Saint-Germain's coach Christophe Galtier has stated that Lionel Messi is not expected to join the team until early January as he is spending time in Argentina following the World Cup. Kylian Mbappé, Neymar Jr. and Achraf Hakimi, who played for their respective national teams at Qatar 2022, could return to the team as long as they are physically and mentally fit. The news is written by Matias Grez from CNN in 2022-12-28 00:00:00. Title: Lionel Messi isn't expected to be back with PSG until early January after World Cup success News:

Outline as Prompt (OP). We also outline the skeleton of original texts by GPT-3.5-turbo and feed the outline into GPT-3 text-davinci-003. The prompt format is as follows:

Prompt for extraction.

"role":	"system",	"con	tent":	"Write	а
hierarchic	al multi-po	oint	outline	for	the
paragraph.	"				
"role": "u	ser", "conte	nt":	{text}		

Example prompt for generation.

960

967

968 969

- 970 971
- 972
- 973

974

957

937

941

943

946

947

951

News Title: There's a shortage of truckers, but TuSimple thinks it has a solution: no driver needed The news is written by Jacopo Prisco, CNN from CNN in 2021-07-15 02:46:59. Outline: I. TuSimple's plan for fully autonomous truck tests A. Reliability of software and hardware needs to improve B. Fully autonomous tests without human safety driver planned by end of year C. Results will determine if company can launch trucks by 2024 D. 7,000 trucks reserved in US alone II. TuSimple's competition Α. Add more details and examples. News

We first remove the HWTs that do not have desired length (i.e., 200-1024 tokens). And we take half of the selected HWTs as references to formulate different prompts mentioned above and feed it into GPT-3 to get MGTs. The MGTs are sampled by Gaussion Distribution of their lengths. To avoid the possible label leakage brought by text length, we directly filter the no-reference HWTs according to the Gaussion Distribution of MGT lengths.

Besides the self-constructed datasets, we also utilize the published GPT-3 dataset TuringBench benchmark (abbraviate as GPT-3 (TB)) (Uchendu et al., 2021) to validate the deceptiveness of GPT-3. The statistics of datasets we use is in Table 6.

Datase	t	Train	Valid	Test	# of tokens
GPT-3(KP)	HWT	446	148	148	427.96 ± 45.49
	MGT	446	148	148	403.88 ± 75.63
GPT-3(SP)	HWT	446	148	148	427.96 ± 45.49
	MGT	446	148	148	415.72 ± 66.54
GPT-3(OP)	HWT	446	148	148	427.96 ± 45.49
	MGT	446	148	148	429.34 ± 78.62
GPT-3(TB)	HWT	5,964	975	1915	236.17 ± 72.96
	MGT	5,507	894	1763	147.29 ± 70.15

Table 6: Statistics of GPT-3 datasets.

We conduct experiments with 3 random seeds and the average results are shown in Table 7. Counterintuitively, even if we elaborate the prompts and eliminate the length difference between MGTs and HWTs, the detection results are still superior, even on outdated baselines like GPT-2. The conclusion might be counterintuitive, but texts generated by the most advanced and popular GPT-3 model are the easiest to detect.

A.3 Implementation Details

This part mentions the implementation details and hyper-parameter settings of all the methods in the experiment. To imitate the situation of low dataresources, we sample 10% texts from the datasets as limited dataset, which will test models together with the complete datasets. And we conduct experiment on 10 different seeds and report the average test accuracy, F1-Score, and standard deviation. 1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1027

1029

1030

1031

1032

1033

1034

1035

1037

1038

1039

1040

1041

We use RoBERTa base model to initialize the embedding of our representation and optimize the model using AdamW (Loshchilov and Hutter, 2017) optimizer with a 0.01 weight decay. We set the initial learning rate to 10^{-5} and the batch size to 8 for all datasets based on experiences.

We utilize packages, namely transformers, pytorch, and allennlp to implement CoCo. And the GPT-3 datasets and ChatGPT case is generated by OpenAI API and websites. We spend \$300 for API costs, including development and final generation costs. We train and do experiments on 8 NVIDIA A100 GPUs on 2 Ubuntu-based servers. The total budget for training 20 epochs, dev, and testing on the GROVER dataset is 2.5 hours. On GPT-2 dataset is 12 hours, and on GPT-3 dataset is 1.5 hours. We will publish our code and dataset recently.

A.4 Effect of Hyper-Parameters

A.4.1 Contrastive Learning Parameters

We evaluate the influence of contrastive learning hyper-parameters α and τ with experiments on different combinations of them. The result is shown in Fig. 5. Considering the discovering that smaller τ leads to better hard negative mining ability (Wang and Liu, 2021), we select α from {0.1, 0.2, ..., 0.9} and τ from {0.1, 0.2, 0.3}. We find that the extreme α value causes the performance degradation and the best hyper-parameter combination is $\alpha, \tau = 0.6, 0.2$. Our analysis is that large α forces the model to concentrate on the instance-level contrast and small α lets class separation objective take control. Both will reduce the generalization performance of the detector on test set.



Figure 5: Effect of parameters α and τ on model performance.

990

991

992

996

997

998

1000

Dataset	GPT-	3 (KP)	GPT-	3 (SP)	GPT-	3 (OP)	GPT-	3 (TB)
Metric	ACC(val/test)	F1(val/test)	ACC (val/test)	F1 (val/test)	ACC (val/test)	F1 (val/test)	ACC (val/test)	F1 (val/test)
GPT2	0.9914/0.9916	0.9916/0.9918	0.9890/0.9893	0.9885/0.9889	0.9925/0.9928	0.9923/0.9924	0.9884/0.5422*	0.9880/0.6335*
RoBERTa	0.9946/0.9950	0.9950/0.9952	0.9935/0.9941	0.9933/0.9937	0.9946/0.9943	0.9942/0.9940	0.9962/0.6406*	0.9960/0.7273*
CoCo	0.9955/0.9950	0.9942/0.9945	0.9938/0.9941	0.9936/0.9940	0.9942/0.9943	0.9942/0.9943	0.9966*	0.9970*

Table 7: Experiment of different detectors on different GPT-3 Dataset. * :The great performance difference between validation set and test set on GPT-3 (TB) are because the test set randomly sample 50% of the words of each article in the dataset (Uchendu et al., 2021). We do not test CoCo on GPT-3 (TB) for the reason that such operation greatly influences the coherence in texts. We provide an example of this in Table 8.

Recent changes to key the unprecedented exclu	011-5(01)
ado cricbuzz.bevan leads scotland's 21-man squad for their first er test match against pakistan in edinburgh icc.chris rogers irres after champions trophy defeat : australian cricketer an- unces international retirement the sun.icc super eight teams bdi ranking results.bahrain host oman on sunday kitply hans hra gold cup gulf today.icc results.new zealand series history ndia v new zealandyazan mohsen qawasma : how bahrain ught construction of the section of the secti	international indexes have resulted in sion of Russian stocks at a "zero" price, in Moscow's already-dismal stock ex- nas made Russia no longer an option for shift to other emerging markets.\n\nThe he in early March, when FTSE Russell he removal of Russian stocks from their try's escalating economic and geopoliti- fter, the Moscow Exchange suspended through the market.\n\nThe possible de- s Western investors further reconsidering isia

Table 8: A comparison example between texts in test set of GPT-3 (TB) and GPT-3 (OP). The GPT-3 (TB) text shows great disorder while GPT-3 (OP) text is neat.

A.4.2 Graph Parameters

'.v lan ev ret no : c vo : i ca

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

We further investigate the effect of max node number and max sentence number on model performance. The result is shown in Fig. 6. We select max node number from {60, 90, 120, 150} and max sentence number from {30, 45, 60, 75}. The detector performs best when max node number is 90 and max sentence number is 45. The experiment results prove that the large node and sentence number are not necessary for the improvement of detection accuracy. We infer that even though setting large node and sentence number includes more entity information, excessive nodes bring noise to the model and impair the distinguishability of coherence feature.



Figure 6: Performance of CoCo with different graph parameters.

A.5 Static Geometric Analysis on Coherence Graph

We have witnessed performance enhancement by applying the graph-based coherence model to the detection model, but how does the coherence graph help detection? In this subsection, we apply static geometric features analysis to coherence graph we construct to evaluate the distinguishable difference between HWTs and MGTs with explanation. In the following discussion, we take the dataset of GROVER into the analysis. Some basic metrics of data and the corresponding graph are shown in Table 9.

Metric	HWT	MGT
Sample Num.	4994	4991
Avg. Num. of Token	463.2	456.0
Avg. Num. of Vertex	43.60	32.37
Avg. Num. of Edge	107.4	65.44

Table 9: Basic metrics of texts and corresponding graphs.

Though HWTs and MGTs have approximately1070the same number of tokens in every text, coher-
ence graph for HWTs has larger scale than MGTs'10711072

1068

1069

1057

1058

1059

Metric	Avg. Degree
HWT	2.980
MGT	2.591

Table 10: Average of degree (whole dataset).

with 34.7% more vertexes and 64.1% more edges,
which shows that HWTs have more complex semantic relation structures than MGTs.

A.5.1 Degree Distribution

Semantically, degree of coherence graph measures the co-occurrence and TF-IDF feature of keywords.
Moreover, degree distribution shows global coherence because high-degree nodes devote to the main topic and low-degree nodes are the extension.



Figure 7: Distribution of average degree of graphs.

As shown in Table 10, the degree of the graph representation of HWTs is **15.0%** larger than MGTs, which shows disparities of MGTs to form coherent interaction between sentences. Fig. 7 measures the distribution of each graph's average nodes' degree, showing that the distribution of HWTs has a longer tail than MGTs.



Figure 8: Distribution of degree with different provenance.

Furthermore, we analyze the distinguishability of degree features when impacted by other factors.

One most considerable influences is the style and genre of different provenance. We chose around 60 articles from The Sun⁹ and Boston¹⁰. Then we use GROVER to mimic their style to generate similar topic news. Fig. 8 shows the degree distribution of HWTs and MGTs of both provenances. 1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

We use Jensen–Shannon divergence to evaluate the similarity of the degree distribution. The JS-divergence of MGTs mimicking The Sun and Boston is **0.029**, while the JS-divergence of MGTs and HWTs in Boston is **0.050**, in The Sun is **0.061**. The apparent gap shows that degree distribution can robustly detect MGTs and HWTs when impacted by provenance differences.

A.5.2 Aggregation

Aggregation is a shared metric for complex networks and linguistics, depicting how closely the whole is organized around its core. We propose two metrics to evaluate the aggregation of graphbased text representation in our coherence model, the size of the largest connected subgraph and the clustering coefficient.

In our representation, not all sentences have entities related to others. Hence the graph is an unconnected one. The average number of nodes in subgraphs of MGTs is **4.49** and of HWTs is **4.84**. We propose that the size of the largest connected subgraph shows the contents which are closely organized around the topic. Moreover, the size of graphs may be an unfair factor, so we use the portion of nodes in the largest connected subgraph to reflect its size. The average portion in HWTs is **0.6725** and in MGTs is **0.6458**. Fig. 9 shows the distribution of the portion of graphs, and HWTs distribute more high-portion ones than MGTs.

The clustering coefficient represents how nodes tend to cluster. For the entities of texts, clustering evaluates how the author narrates around the central theme. The larger the clustering coefficient is, the tighter the semantic structure is. The average cluster coefficient of the graphs of HWTs is **0.2213** and of MGTs is **0.1983**, HWTs is **11.6%** better than MGTs. Fig. 10 shows the distribution.

A.5.3 Core & Degeneracy

The degeneracy of a graph is a measure of how sparse it is, and the *k*-core is the subgraph corresponding to its significance in the graph. We

1076

1079

1080

1081

1082

1083

1084

1085

1086

1087

⁹https://www.thesun.co.uk/

¹⁰https://www.boston.com/



Figure 9: Portion of the largest connected subgraph.



Figure 10: Distribution of clustering coefficient.

1138 propose that, in our graph representation, the degeneracy process of graphs equals summarizing 1139 texts semantically. The maximum of core-number 1140 shows the complexity of hierarchical structure in 1141 texts. Furthermore, the distribution of the core-1142 1143 number reflects the overall sparse and is a graphperspective N-gram module. Based on experiments, 1144 the average core-number of HWTs is 5.772 while 1145 MGTs with 4.458. HWTs are 29.5% ahead. Fig. 11 1146 is the distribution of the core-number. 1147

A.5.4 Entropy

Entropy is a scientific concept to measure a state of disorder, randomness, or uncertainty. The wellknown Shannon entropy is the core of the information theory, measuring the self-information content. For the graph data, network structure entropy defined as the following can examine the information amount of the graph structure.

$$Entropy = -\sum_{i=1}^{N} I_i \ln I_i = -\sum_{i=1}^{N} \frac{k_i}{\sum_{j=1}^{N} k_j} \ln(\frac{k_i}{\sum_{j=1}^{N} k_j}),$$
(9)



Figure 11: Core-number of nodes in graphs



Figure 12: Structure entropy of graphs

where I_i is the information content represented by the degree distribution, N is the number of nodes, and k_i is the degree of the *i*-th node. 1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

Global coherence, from our perspective, equals refining more information inside the semantic structure of the whole text, which matches to structure entropy of our graph representation. From our experiments, the structure entropy of HWTs (2.263) is **6.80%** larger than MGTs (2.119), which means HWTs obtain more structured information because their semantic information is globally organized. We show the network structure entropy distribution in Fig. 12.

A.6 Exploration on Imbalanced Data

Imbalanced distribution in data is another crucial 1171 limitation in the task of MGTs detection, which is 1172 similar to the low resource limitation. It is imag-1173 inable that, with the development of generation 1174 technology, MGTs will overwhelmingly dominate 1175 low-quality articles since they are easier and faster 1176 to generate than human writing. The detection 1177 model will face training resources with MGTs as 1178

1156

1148

1149

1150

1151 1152

1153

1154

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1202

1203

1204

1205

1206

1207

1208

1210

1211

1212

1213

1214

1215

1216

1217

1218

the main part and HWTs as the small part. We test the current models in the imbalanced limitation and find the dramatic decline in accuracy when the ratio of HWTs is less than 30%, as shown in the Fig. 13. The test is based on the 10% GROVER dataset.



Figure 13: Model comparison results on DL dataset with 9 different human-generated text portions.

All models show poor performance at low HWTs ratios. With a percentage of HWTs of 0.1 (only 100 HWTs in the training set in this case), most of the models have an accuracy below 50%, which performance is close to random and reflects intolerance for extreme cases. Besides, we find that a high proportion of HWTs also cause a decrease in F1 score to some extent.

A.7 Related Work: Graph-based Text Representation

Graph can represent text structure and innerrelation (Minsky, 1982). Based on different organizing methods, graphs can reflect static statistical (*e.g.*, co-occurrence (Cancho and Solé, 2001), collocation (Bordag et al., 2003)), dynamically statistical (*e.g.*, evolution (Dorogovtsev and Mendes, 2001)), lexical (Widdows and Dorow, 2002), orthographic (Choudhury et al., 2007), cognitive (*e.g.*, conception (Motter et al., 2002)), syntactic (Cancho et al., 2004; Ferrer Cancho et al., 2007), semantic (Steyvers and Tenenbaum, 2005; Sigman and Cecchi, 2002; Kozareva et al., 2008) relations. Hassan et al. (2017) propose an overall survey about graph-based text representation.

Graph-of Words (GoW) Model (Turney, 2002; Mihalcea and Tarau, 2004) is a type graph representation method in which each document is represented by a graph, whose nodes correspond to terms and edges capture co-occurrence relationships between terms. Using GoW, keywords can be extracted by retaining the document graph (Turney, 2002). Thus, graph representation is sensible to apply in tasks like information retrieval (Blanco and Lioma, 2011), categorization (Malliaros and Skianis, 2015) and sentiment classification tasks (Huang and Carley, 2019; Hou et al., 2021). Most models enhance classification or detection performance by combining graph representation with neural networks. Text-GCN (Yao et al., 2019) first builds a single large graph for the whole corpus, followed by Tensor-GCN (Liu et al., 2020) with tensor representation. Also, the relation between words varies, and should be treated as different edges. CoCo matches keywords PLM embedding to nodes and sentence representation, considers dealing inner- and inter-sentence relation differently in GCN, and merges the structure graph and flat sequence representation to predict accurately.

1219

1220

1221

1222

1224

1225

1226

1227

1228

1230

1231

1232

A.8 Pseudocode of COCO

Algorithm 1 Algorithm of CoCo

- **Input:** Input X, consisting of documents D and corresponding coherence graph G, hyperparameters such as the size of dynamic memory bank M and batch size S, labels Y
- **Output:** A learned model COCO, consisting of key encoder f_k with parameters θ_k , query encoder f_q with parameters θ_q , classifier f_c with parameters θ_c
- 1: Initialize $\theta_k = \theta_q, \theta_c$
- 2: Initialize dynamic memory bank with $f_k(x_1, x_2...x_M)$, where x_i is randomly sampled from X.
- 3: Freeze θ_k
- 4: $epoch \leftarrow 0$
- 5: while $epoch \leq epoch_{max}$ do
- $6: \quad n \leftarrow 0$
- 7: while $n \leq n_{\max} \operatorname{do}$
- 8: Randomly select batch $\boldsymbol{b}_k, \boldsymbol{b}_q$
- 9: $\boldsymbol{D}_q = f_q(\boldsymbol{b}_q), \, \boldsymbol{D}_k = f_k(\boldsymbol{b}_k)$
- 10: $\widehat{p} = f_c(\boldsymbol{D}_q)$
- 11: Calculate $\hat{\mathcal{L}}_{ICL}$ with equation 5, calculate \mathcal{L}_{CE} with equation 6, calculate \mathcal{L}_{total} with equation 7
- 12: Backward on \mathcal{L}_{total} and update θ_q , θ_c based on AdamW gradient descent with an adjustable learning rate
- 13: Momentum update θ_k with equation 8
- 14: Update dynamic memory bank queue with $enqueue(queue, D_k), dequeue(queue)$

15: $k \leftarrow k+1$

- 16: end while
- 17: **if** Early stopping **then**
- 18: break
- 19: else
- 20: $epoch \leftarrow epoch + 1$
- 21: end if
- 22: end while
- 23: return A trained model CoCo